# MACHINE LEARNING AND ARTIFICIAL INTELLIGENCE ON AWS

Lucian Revnic

# AGENDA

Introduction

Infrastructure
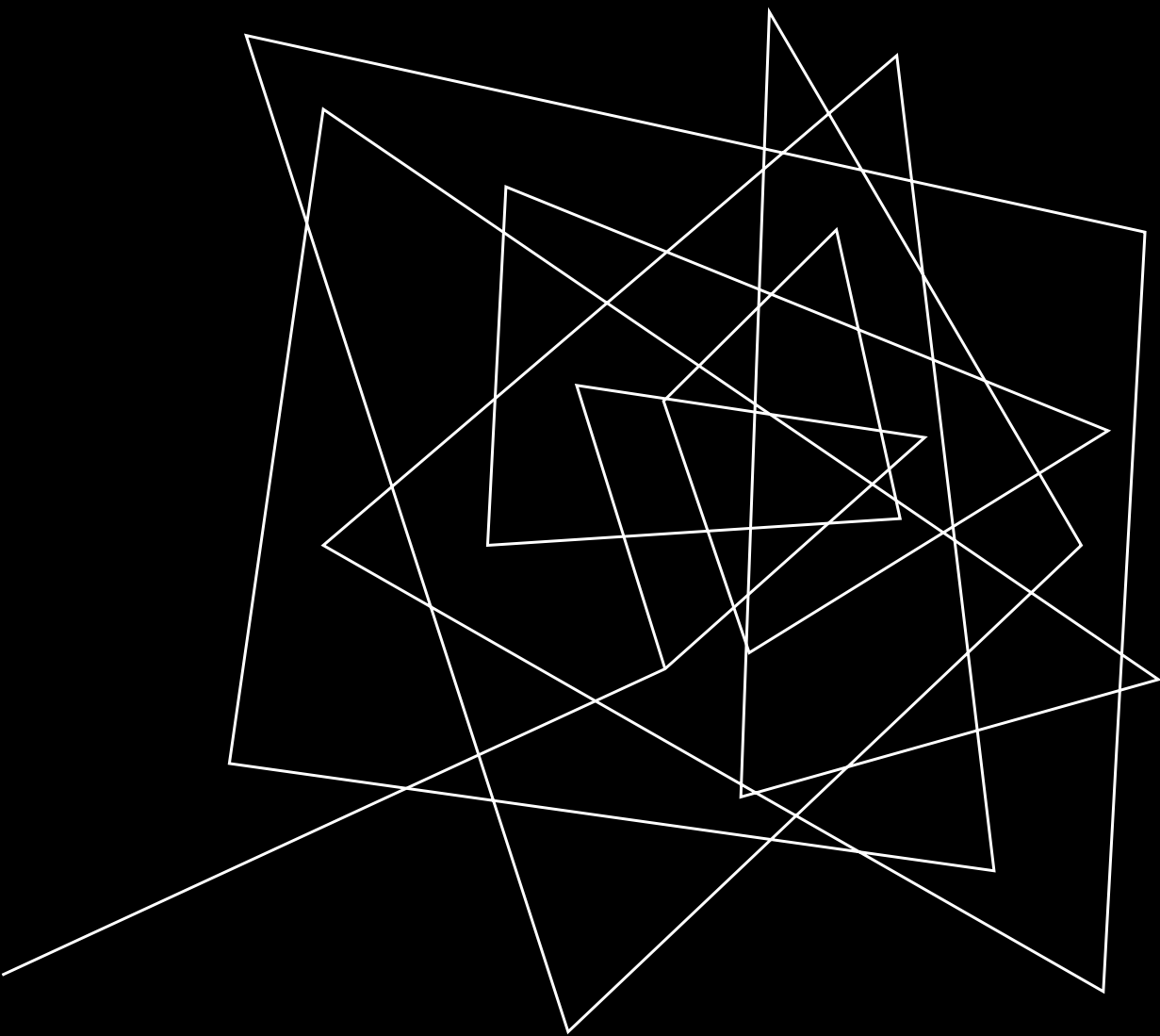
Tools

Applications and Services

INFRASTRUCTURE FOR TRAINING AND INFERENCE
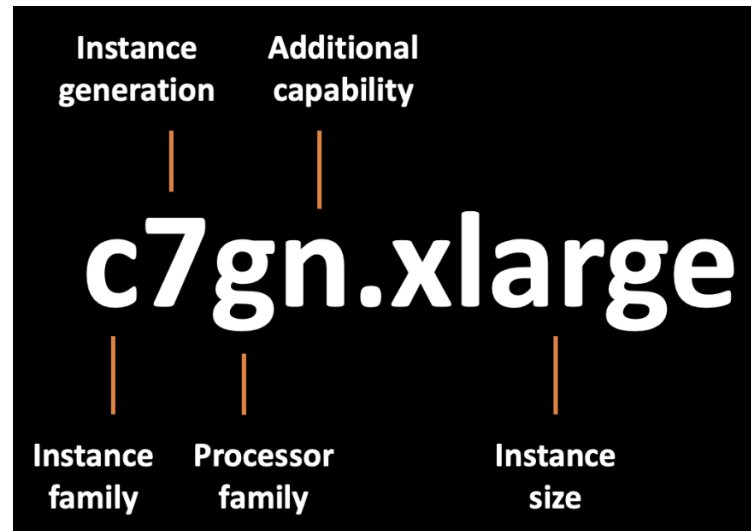
BUILD AND DEVELOPMENT TOOLS

APPLICATION SERVICES

# INFRASTRUCTURE

# AWS INSTANCE TYPES – QUICK REFRESH

## Instance families

- **C** – Compute optimized
- **D** – Dense storage
- **F** – FPGA
- **G** – Graphics intensive
- **Hpc** – High performance computing
- **I** – Storage optimized
- **Im** – Storage optimized with a one to four ratio of vCPU to memory
- **Is** – Storage optimized with a one to six ratio of vCPU to memory
- **Inf** – AWS Inferentia
- **M** – General purpose
- **Mac** – macOS
- **P** – GPU accelerated
- **R** – Memory optimized
- **T** – Burstable performance
- **Trn** – AWS Trainium
- **U** – High memory
- **VT** – Video transcoding
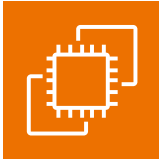- **X** – Memory intensive

## Processor families

- **a** – AMD processors
- **g** – AWS Graviton processors
- **i** – Intel processors

## Additional capabilities

- **d** – Instance store volumes
- **n** – Network and EBS optimized
- **e** – Extra storage or memory
- **z** – High performance
- **q** – Qualcomm inference accelerators
- **flex** – Flex instance

# INFRASTRUCTURE FOR TRAINING AND INFERENCE

**Amazon Elastic Compute Cloud (Amazon EC2)**

**AWS Deep Learning AMIs**

**AWS Deep Learning Containers**

**EC2 P5**

**GPU base instances / GPU accelerated**
8 NVIDIA® V100 Tensor Core GPUs
Up to 3200 Gbps net.
Deployed in EC2 **Ultra Clusters**
$$$$$

**EC2 C5**
Cheaper that GPU (P*, G*) instances
Intel or AMD processors
100 Gbps of net.
$

**EC2 G5**

**Graphics intensive**
NVIDIA A10G
Most cost-effective GPU instances inference & graphics
$$

**EC2 Inf2**

Reduce cost on **ML prod** deployments
Up to 12 AWS Inferentia2 accelerators
Up to 100 Gbps networking
$$

**EC2 Tr1**
Deep Learning Training
AWS Trainium chips
Intel Xel Scalable 3rd generation
Up to 16  AWS Trainium accelerators
Up to 1600 Gbps networking
Deployed in **EC2 Ultra Clusters**
$$$

x86 or AWS Graviton
Processor :GPU vs CPU vs Inferetia vs Hababa
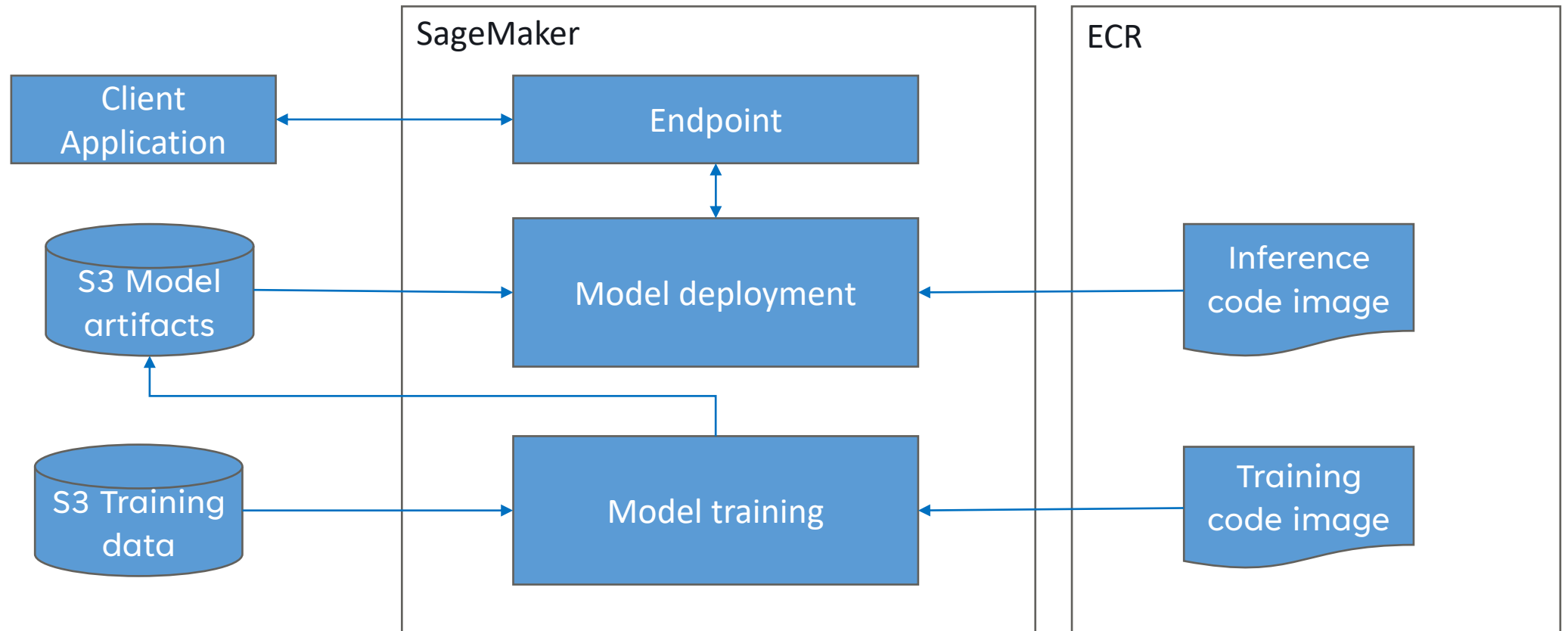SDK: CUDA vs AWS Neuron vs SynapsesAI
OS: Amazon Linux vs Ubuntu

60 DLCs for 5x frameworks (TensorFlow, PyTorch, MXNET, AutoGluon, Hugging Face)
EC2, ECS, EKS, Fargate, Lambda
General containers and AWS Neuron containers
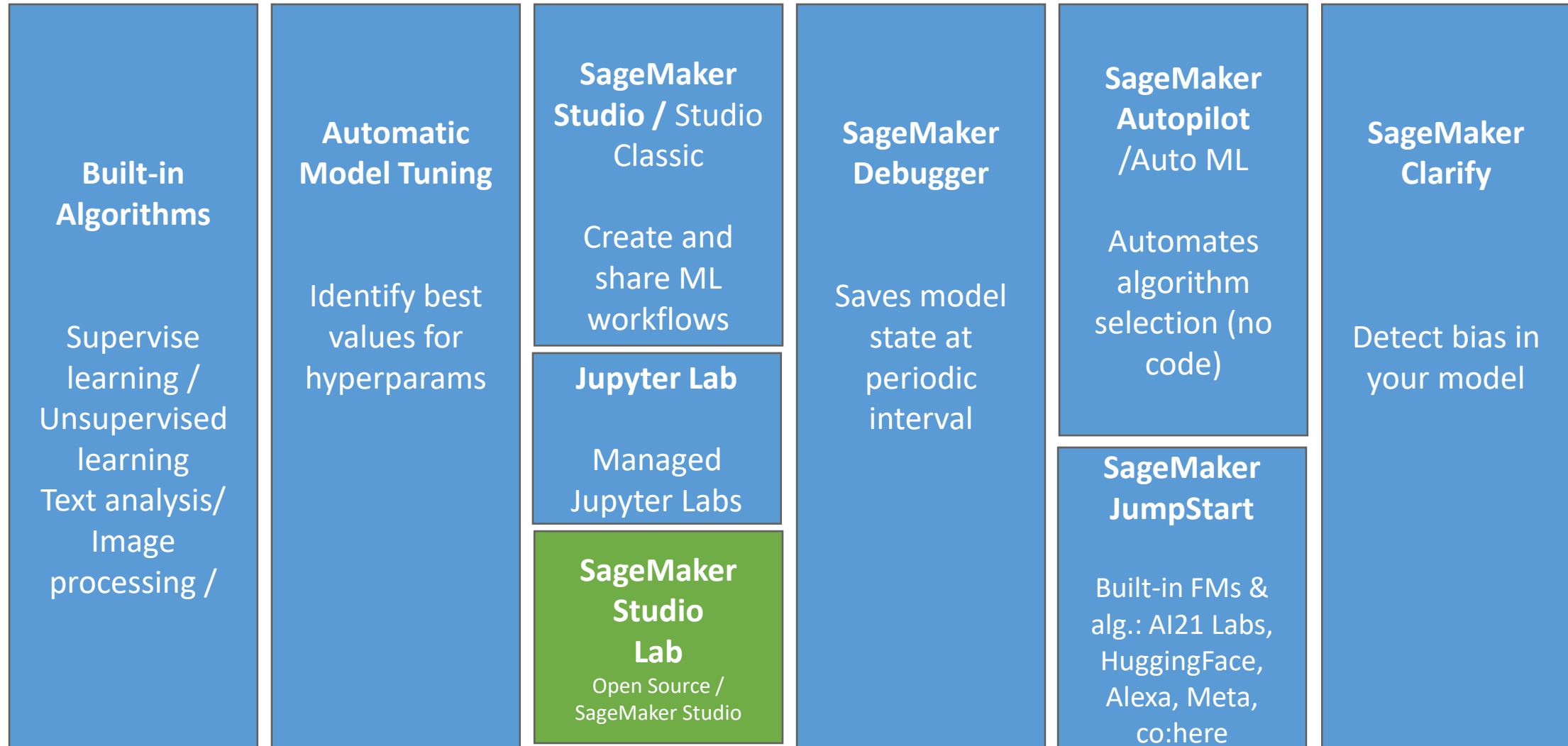For AWS or on-premises (eg ECS Anywhere, EKS Anywhere)

# BUILD ML MODELS SERVICES AND TOOLS

# AMAZON SAGEMAKER (1)

## Fully managed machine learning (ML) service to build, train, and deploy ML models



https://www.sundog-education.com/

# AMAZON SAGEMAKER (2) – BUILD & TRAIN

**Built-in Algorithms**

Supervise learning / Unsupervised learning Text analysis/ Image processing /

**Automatic Model Tuning**

Identify best values for hyperparams

**SageMaker Studio /** Studio Classic

Create and share ML workflows

**Jupyter Lab**

Managed Jupyter Labs

**SageMaker Studio Lab**

Open Source / SageMaker Studio

**SageMaker Debugger**

Saves model state at periodic interval

**SageMaker Autopilot** /Auto ML

Automates algorithm selection (no code)

**SageMaker JumpStart**

Built-in FMs & alg.: AI21 Labs, HuggingFace, Alexa, Meta, co:here

**SageMaker Clarify**

Detect bias in your model

# AMAZON SAGEMAKER (3) – DEPLOY BASED ON USE CASES



Source: https://www.softwebsolutions.com/resources/guide-to-amazon-sagemaker.html

## Real Time Inference
Persistent, real-time endpoints

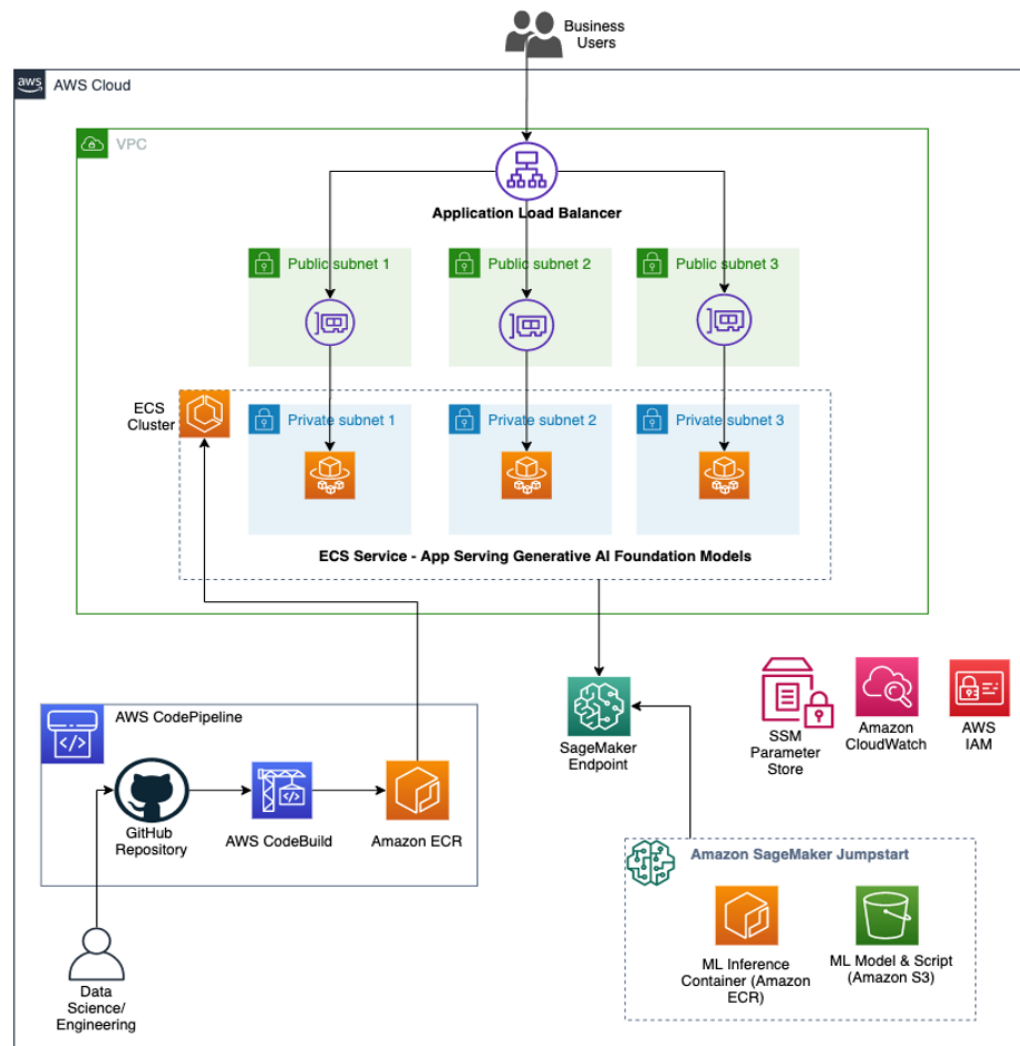## Serverless Inference
Workloads that have idle periods

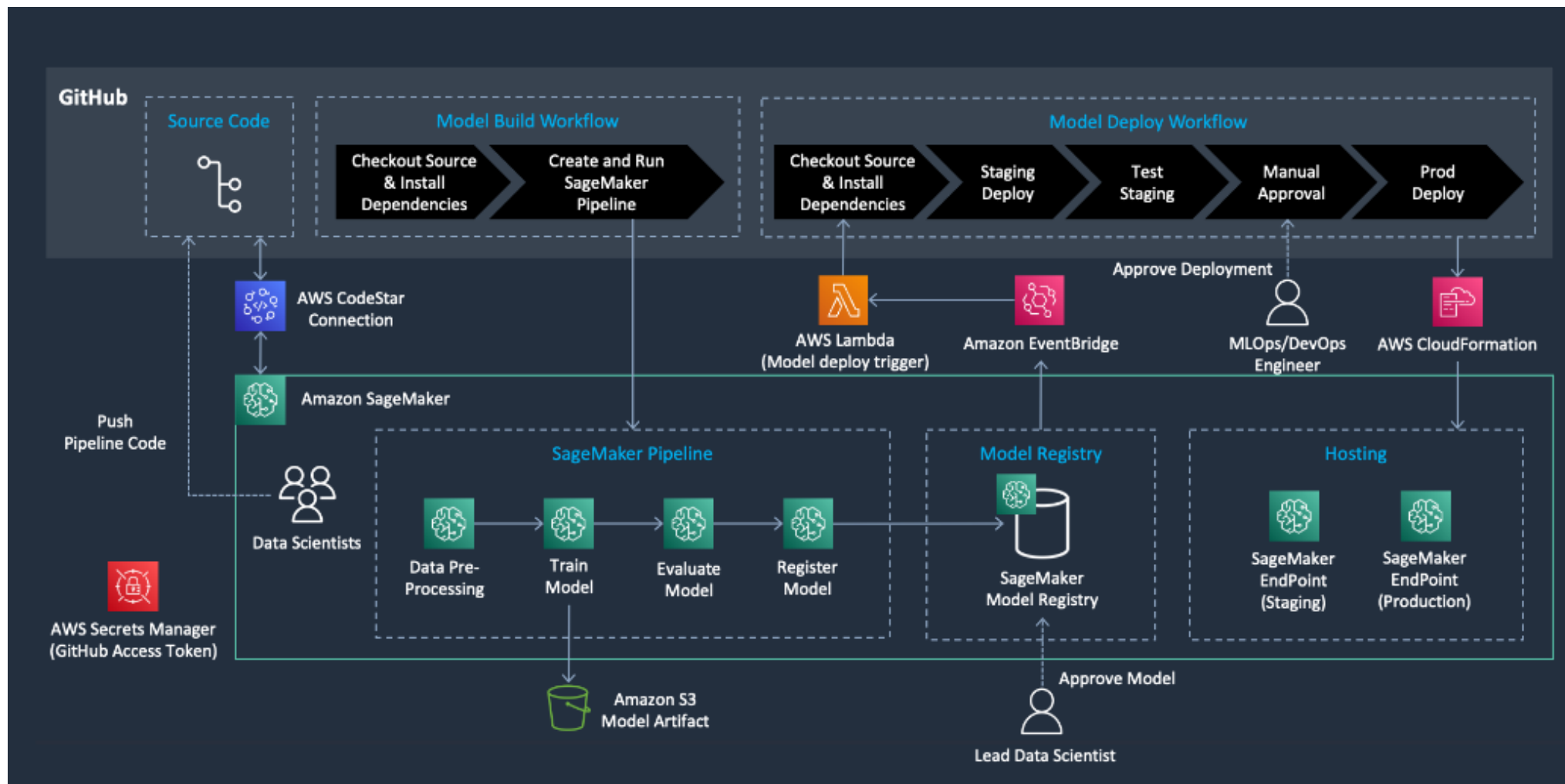## Asynchronous Inference
Large payloads (eg up to 1GB)

## Batch Transform
Prediction on entire data set
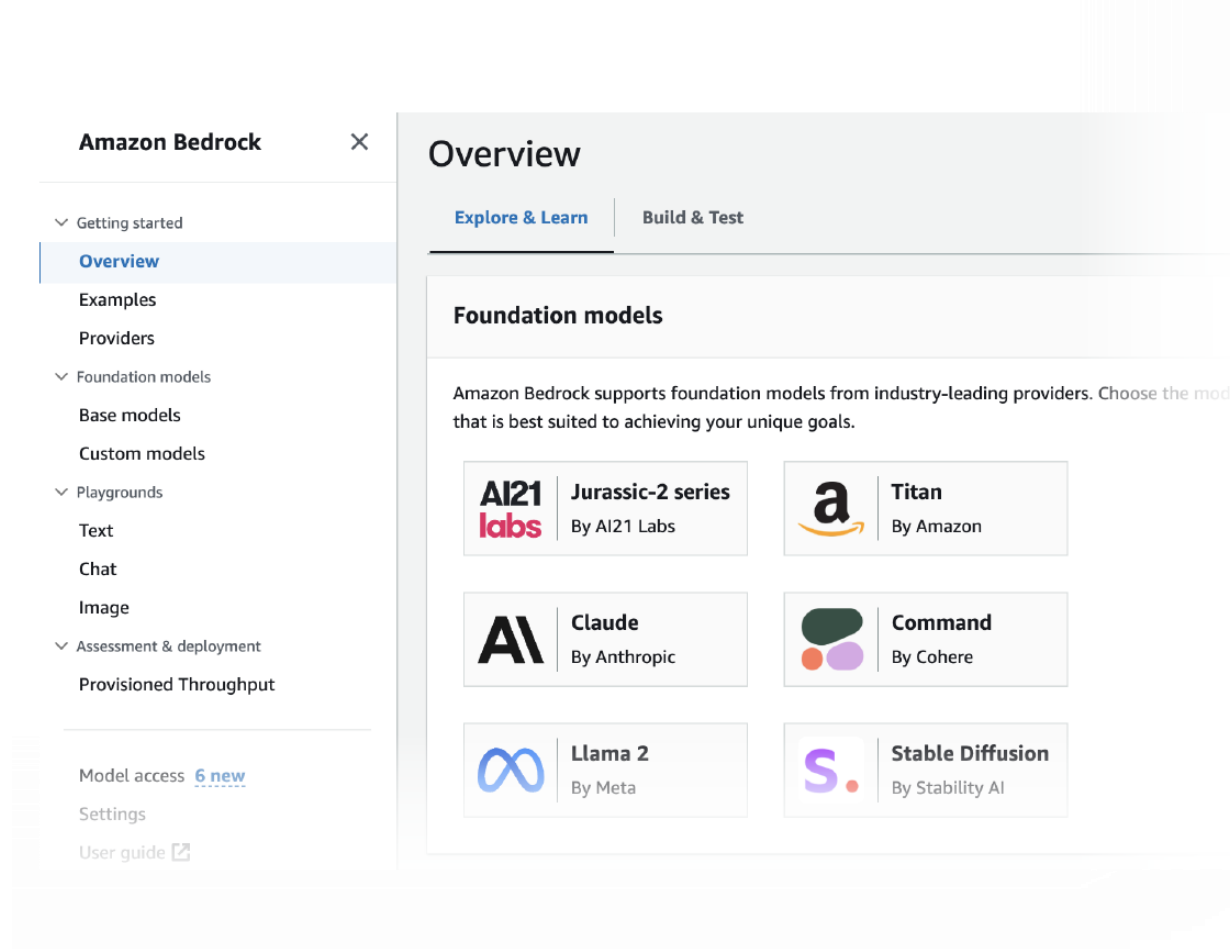
# AMAZON SAGEMAKER – EXAMPLES (GENERATIVE AI APP)



https://aws.amazon.com/blogs/containers/build-generative-ai-apps-on-amazon-ecs-for-sagemaker-jumpstart/

# AMAZON SAGEMAKER – EXAMPLES (MLOPS PIPELINE)

# AMAZON BEDROCK

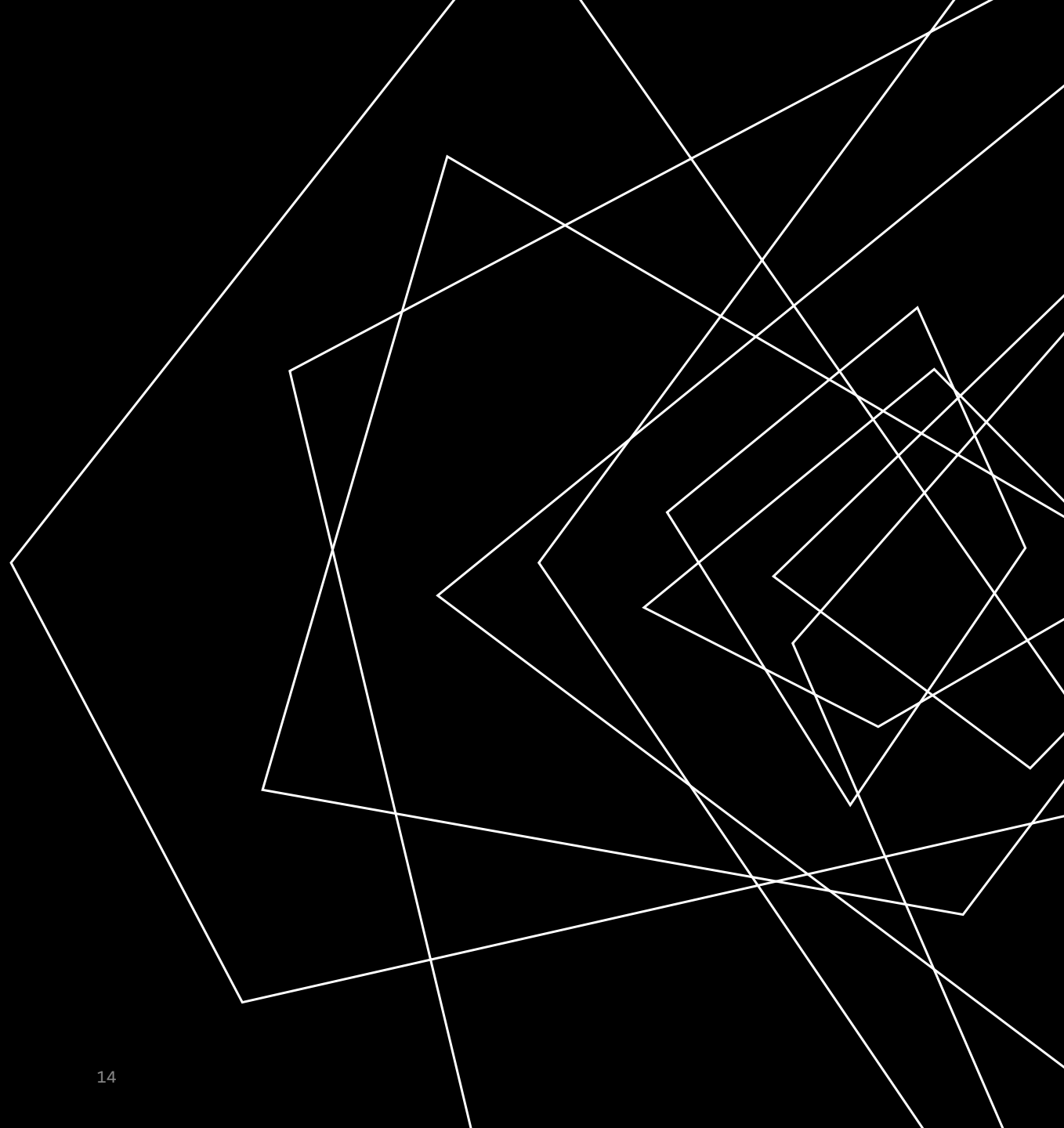The easiest way to build and scale generative AI applications with foundation models



Platform for FMs as it enables you to select your preferred model

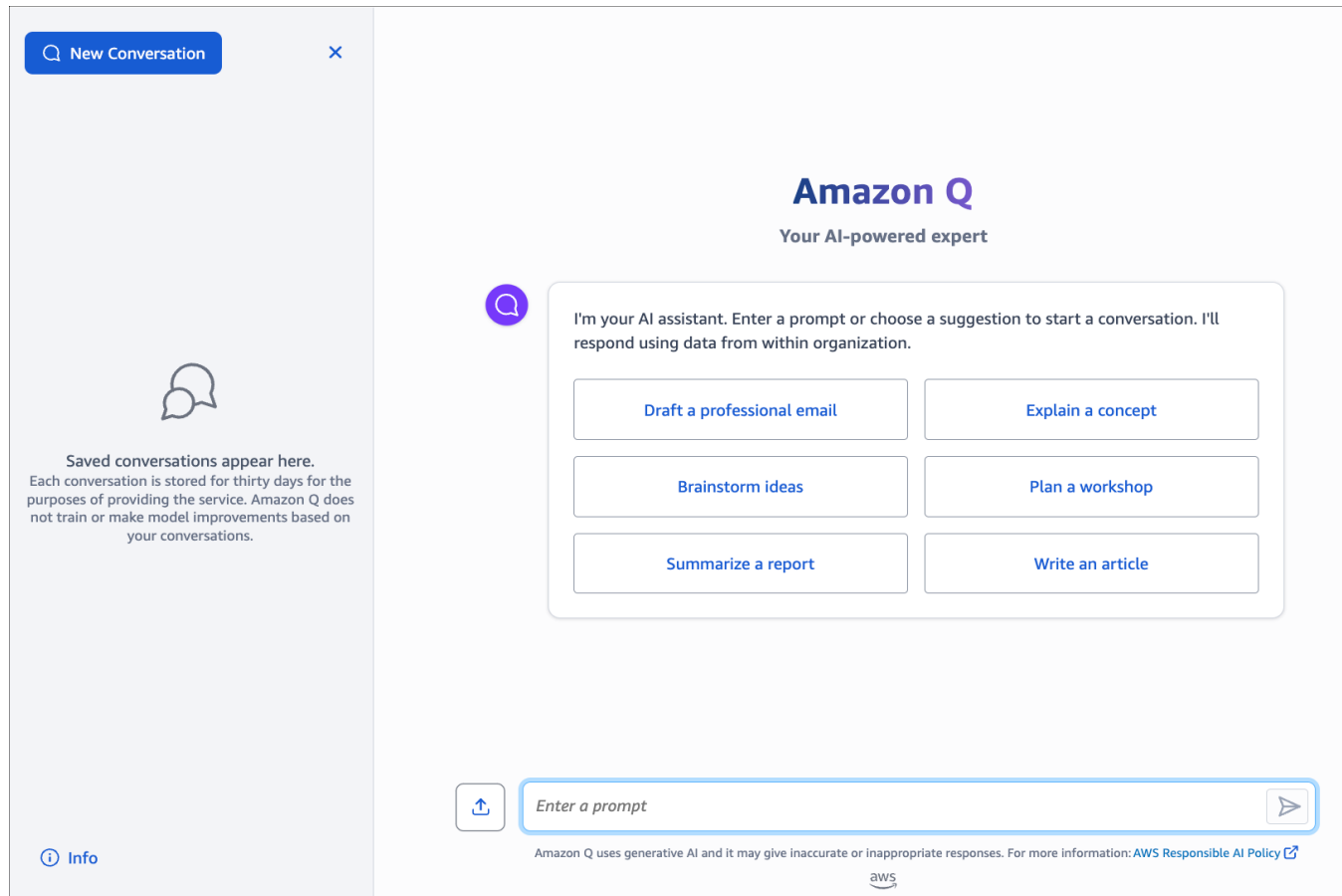Single API for inference

Playgrounds

# AI APPLICATION SERVICES
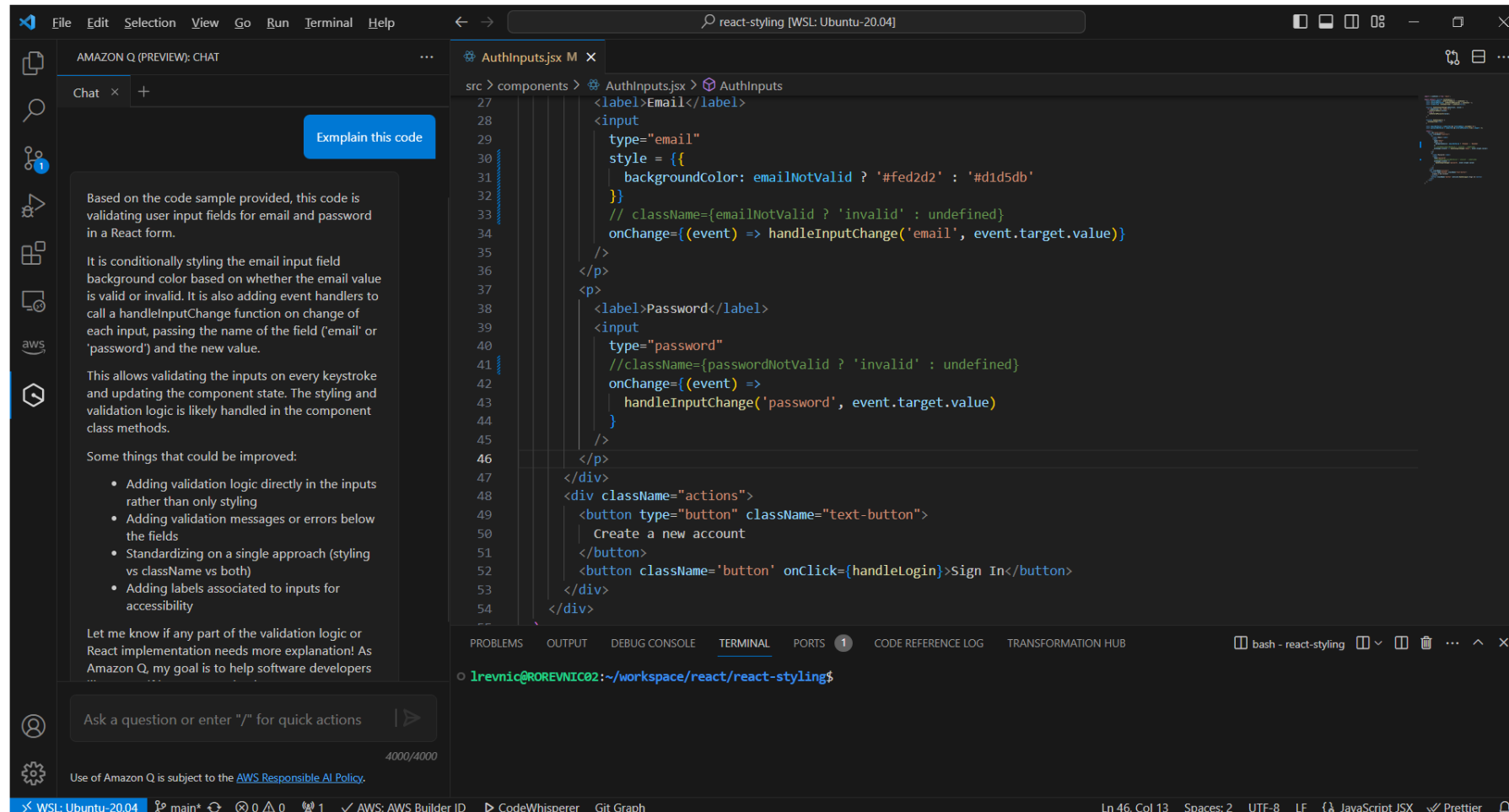
# AMAZON Q

## The new generative AI–powered assistant



Connects to your data such as Jira, GitHub, ServiceNow

Integrated with AWS Console: "Why is my instance not accessible?"

Integrated with AWS QuickSight: "Build me a strategy to improve sales"

# AMAZON Q

**NEW!**



Visual Studio Code Integration via AWS Toolkit



Amazon Q
Your generative AI-powered
assistant tailored for work

# AMAZON CODEWHISPERER



Code generation

Security scanning

IaC code generation (HashiCorp Terraform, AWS CloudFormation, CDK)

**NEW!**

# OTHER ML AND AI SERVICES

**Amazon Comprehend**
*NLP parsing and understanding*

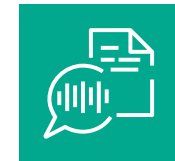**Amazon Elastic Inference**
*GPU Management - Deprecated*

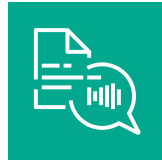**Amazon Forecast**
*Forecasting*

**Amazon Lex**
*Conversational interfaces*

**Amazon Transcribe**
*Speech to Text*

**Amazon Personalize**
*Recommendations*

**Amazon Polly**
*Text to Speech*

**Amazon Rekognition**
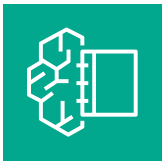*Image recognition & video analysis*

**Amazon Textract**
*Document understanding / OCR*

**Amazon Translate**
*Machine translation*

**Amazon SageMaker Studio Lab**

**Amazon Comprehend Medical**
*Medical text and documents*

**AWS HealthOmics**
(successor to Amazon Omics)

**AWS Neuron**
*Deep Learning SDK*

**Amazon CodeGuru**
*Detect and FIX code vulnerabilities*

**Amazon DevOps Guru**
*Apps anomaly detection*

**AWS HealthLake**
(successor to Amazon HealthLake)

**AWS Panorama**
*Computer Vision*

**AWS HealthImaging**
*Medical Images in the Cloud*

**AWS HealthScribe**
*Clinical document generation*

NEW!

AWS HealthScribe

# RESOURCES

AI/ML AWS Workshops

https://workshops.aws/categories/AI%2FML



AWS Trainings – FREE!

https://skillbuilder.aws/



AWS Machine Learning University

https://aws.amazon.com/machine-learning/mlu/

# THANK YOU

## Lucian Revnic

lrevnic@gmail.com

## Transsylvania Cloud

https://www.meetup.com/TransylvaniaCloud/